

# Benchmark on Automatic 6-month-old Infant Brain Segmentation Algorithms: The iSeg-2017 Challenge

Li Wang<sup>\*,\*</sup>, Senior Member, IEEE, Dong Nie<sup>‡</sup>, Guannan Li<sup>‡</sup>, Élodie Puybureau<sup>‡</sup>, Jose Dolz<sup>‡</sup>, Qian Zhang<sup>‡</sup>, Fan Wang<sup>‡</sup>, Jing Xia<sup>‡</sup>, Zhengwang Wu<sup>‡</sup>, Jiawei Chen, Toan Duc Bui, Jitae Shin, Guodong Zeng, Guoyan Zheng, Vladimir S. Fonov, Andrew Doyle, Yongchao Xu, Pim Moeskops, Josien P.W. Pluim, Christian Desrosiers, Ismail Ben Ayed, Gerard Sanroma, Oualid M. Benkarim, Adrià Casamitjana, Verónica Vilaplana, Gang Li, Member, IEEE, Weili Lin, and Dinggang Shen<sup>\*,</sup> Fellow, IEEE

**Abstract**—Accurate segmentation of infant brain magnetic resonance (MR) images into white matter (WM), gray matter (GM), and cerebrospinal fluid (CSF) is an indispensable foundation for early studying of brain growth patterns and morphological changes in neurodevelopmental disorders. Nevertheless, in the isointense phase (approximately 6-9 months of age), due to inherent myelination and maturation process, WM and GM exhibit similar levels of intensity in both T1-weighted (T1w) and T2-weighted (T2w) MR images, making tissue segmentation very challenging. Despite many efforts devoted to brain segmentation, only few studies have focused on the segmentation of 6-month infant brain images. With the idea of boosting methodological development in the community, iSeg-2017 challenge (<http://iseg2017.web.unc.edu>) provides a set of 6-month infant subjects with manual labels for training and testing the participating methods. Among the 21 automatic segmentation methods participating in iSeg-2017, we review the 8 top-ranked teams, in terms of Dice ratio, modified Hausdorff distance and average surface distance, and introduce their

pipelines, implementations, and source codes. We further discussed limitations and possible future directions. We hope the dataset in iSeg-2017 and this review article could provide insights into methodological development for the community.

**Index Terms**—Infant, brain, segmentation, isointense phase, challenge

## I. INTRODUCTION

THE first year of life is the most dynamic phase of the postnatal human brain development, along with rapid tissue growth and development of a wide range of cognitive and motor functions. The increasing availability of non-invasive infant brain multimodal magnetic resonance images (MRI), e.g., T1-weighted (T1w) and T2-weighted (T2w) images, provides unprecedented opportunities for accurate and reliable charting of dynamic early brain developmental trajectories in understanding normative and aberrant growth. For example, the Baby Connectome Project<sup>1</sup> (BCP) [1] is acquiring and releasing both cross-sectional and longitudinal high-resolution multimodal MRI data from 500 typically-developing children from birth to 5 years of age. The Developing Human Connectome Project<sup>2</sup> (dHCP) in the UK is releasing MRI data from 1500 subjects acquired from 20 to 44 weeks post-conceptional age. These large-scale datasets will undoubtedly greatly increase our limited knowledge on normal early brain development, and provide important insights into the origins and abnormal developmental trajectories of neurodevelopmental disorders, such as autism, schizophrenia, bipolar disorder, and attention-deficit/hyperactivity disorder.

One fundamentally important step in studying the normal and abnormal early brain development is accurate segmentation of infant brain MR images into different regions of interest (ROIs), e.g., white matter (WM), gray matter (GM), and cerebrospinal fluid (CSF). There are three distinct phases in the first-year brain MRI, as shown in Fig. 1. During the infantile phase ( $\leq 5$  months), GM shows higher signal intensity than WM in T1w images. The isointense phase (6-9 months) corresponds to the myelination and maturation process of the brain, yielding to an increase of the intensity of WM in T1w images and thus a low signal differentiation between GM and WM (which is also the case for T2w

This work was supported in part by National Institutes of Health grants MH109773, MH117943, MH100217, MH070890, EB006733, EB008374, EB009634, AG041721, AG042599, MH088520, MH108914, and MH107815.

<sup>\*,\*</sup> L. Wang, <sup>‡</sup> D. Nie, <sup>‡</sup> G. N. Li, Q. Zhang, F. Wang, J. Xia, Z. Wu, J. Chen, G. Li and W. Lin are with the Department of Radiology and Biomedical Research Imaging Center, UNC-Chapel Hill, NC, 27599 USA. (e-mail: [li\\_wang@med.unc.edu](mailto:li_wang@med.unc.edu))

T. Bui and J. Shin are Media System Lab., School of Electronic and Electrical Eng., Sungkyunkwan University (SKKU), Korea

G. Zeng and G. Zheng are with Information Processing in Medical Intervention Lab., University of Bern, Switzerland.

V. Fonov is with NeuroImaging and Surgical Technologies Lab, Montreal Neurological Institute, McGill University, Montreal, Quebec, Canada.

A. Doyle is with McGill Centre for Integrative Neuroscience, Montreal Neurological Institute, McGill University, Montreal, Quebec, Canada.

Y. Xu and <sup>‡</sup> É. Puybureau are with EPITA Research and Development Laboratory (LRDE), Le Kremlin-Bicêtre, France.

P. Moeskops and J. Pluim are with Medical Image Analysis Group, Department of Biomedical Engineering, Eindhoven University of Technology, Eindhoven, Netherlands.

<sup>‡</sup> J. Dolz, C. Desrosiers, and I. Ayed are with Laboratory for Imagery, Vision and Artificial Intelligence (LIVIA), Ecole de technologie supérieure, Montreal, Canada.

G. Sanroma and O. Benkarim are with Simulation, Imaging and Modelling for Biomedical Systems (SIMBIOsys), Universitat Pompeu Fabra, Spain.

A. Casamitjana and V. Vilaplana are with Universitat Politècnica de Catalunya, Barcelona Tech, Spain.

<sup>‡</sup> D. Shen is with the Department of Radiology and Biomedical Research Imaging Center, UNC-Chapel Hill, NC, USA; and Department of Brain and Cognitive Engineering, Korea University, Seoul 02841, Republic of Korea. (e-mail: [dgshen@med.unc.edu](mailto:dgshen@med.unc.edu))

<sup>‡</sup> Co-first authors

<sup>\*</sup> Co-corresponding authors

<sup>1</sup> <http://babyconnectomeproject.org>

<sup>2</sup> <http://www.developingconnectome.org>

images). The last phase is the early adult-like phase ( $>9$  months), where GM intensity is much lower than that of WM in T1w images, similar to the pattern of tissue contrast in the adult MR images. The corresponding tissue intensity distributions of three phases are shown in the third row of Fig. 1, from which we can see the relative good contrast for the infantile and early adult-like phases. However, in the isointense phase, the intensity distributions of voxels in GM and WM are largely overlapping (especially in the cortical regions), thus leading to the lowest tissue contrast and creating the main challenge for tissue segmentation, in comparison to images at other phases of brain development. Also, the appearance of exact isointense contrast varies across different brain regions due to nonlinear brain development [2]. These patterns, along with various factors, such as motion artifacts or severe partial volume effect due to smaller brain size and ongoing white matter myelination, make automatic segmentation of isointense infant brain MRI a highly challenging task, thus causing that existing computational tools typically developed for processing and analyzing adult brain MRI, e.g., SPM, FSL, BrainSuite, CIVET, FreeSurfer and HCP pipeline, often perform poorly on infant brain MRI [3].

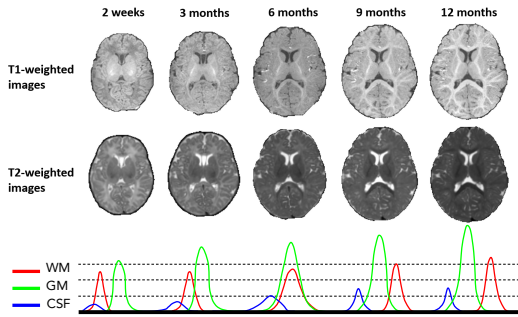


Fig. 1. T1- and T2-weighted MR images of an infant longitudinally scanned at 2 weeks, 3, 6, 9 and 12 months of age. At around 6 months of age, MR images show the lowest tissue contrast, implying the most significant challenge for tissue segmentation. The corresponding tissue intensity distributions from T1w MR images are shown in the bottom row, which indicates high overlap of WM and GM intensities in the isointense phase.

We have witnessed the spread and rise in popularity of Grand Challenges in the medical imaging community during the last years (e.g., NeoBrainS12<sup>3</sup> [4], MRBrains<sup>4</sup> [5], ISLES<sup>5</sup> [6], and BRATS<sup>6</sup> [7]). These challenges have allowed development of public benchmarks that serve as fair and up-to-date comparisons for the methods proposed by colleagues around the world. For example, the MICCAI challenge on neonatal MRI segmentation (NeoBrainS12<sup>3</sup>) and the MICCAI challenge on adult MRI segmentation (MRBrains<sup>4</sup>) mainly focused on the infantile and adult-like phases, respectively, rather than the challenging isointense phase. To date, only a few studies focused on the segmentation of 6-month infant brain image [4, 8-10]. In iSeg-2017<sup>7</sup> challenge, researchers were invited to participate

with their automatic algorithms to segment WM, GM and CSF on isointense (6-month) infant brain MR scans, which remains scarce in the field. At the time of writing this paper, 21 teams had submitted their results on the iSeg-2017 website. In this paper, we focus only on those methods that were ranked among the 8 top teams in terms of Dice Coefficient (DICE), modified Hausdorff distance (HD95) and Average Surface Distance (ASD). In the next section, we introduce the cohort employed for this challenge. Then, in Section III, the metrics used to evaluate the performance of the proposed methods are detailed. Section IV provides a complete description of the top-ranked methods selected for this review. Section V discusses their performance, limitations and possible future directions.

## II. DATA

Selected MR scans for training and testing were randomly chosen from the pilot study of Baby Connectome Project (BCP, <http://babyconnectomeproject.org>). All infants were term born (GA  $40 \pm 1$  weeks) without any pathology. At the time of scanning, average age is  $6.0 \pm 0.5$  months old. MR scans were acquired on a Siemens head-only 3T scanners with a circular polarized head coil. During the scan, infants were asleep, unsedated, fitted with ear protection, and their heads were secured in a vacuum-fixation device.

- T1-weighted MR images were acquired with 144 sagittal slices using parameters: TR/TE = 1900/4.38 ms, flip angle =  $7^\circ$ , resolution =  $1 \times 1 \times 1$  mm<sup>3</sup>;
- T2-weighted MR images were obtained with 64 axial slices: TR/TE = 7380/119 ms, flip angle =  $150^\circ$ , resolution =  $1.25 \times 1.25 \times 1.95$  mm<sup>3</sup>.

For image preprocessing, T2w images were rigidly aligned onto their corresponding T1w images. All images were resampled into an isotropic  $1 \times 1 \times 1$  mm<sup>3</sup> resolution. Next, standard image preprocessing steps were performed before manual segmentation, including skull stripping [11], intensity inhomogeneity correction [12], and manual removal of the cerebellum and brain stem by experts.

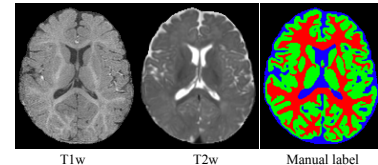


Fig. 2. T1w and T2w MR images of an infant subject scanned at 6 months of age (isointense phase), provided by iSeg-2017. From left to right: T1w MR image, T2w MR image, and manual label image.

To generate reliable manual segmentation, we first took advantage of the follow-up 24-month scans of the same subjects, with high tissue contrast, to generate an initial automatic segmentation for 6-month scans by using a publicly available software iBEAT ([www.nitrc.org/projects/ibeat/](http://www.nitrc.org/projects/ibeat/)). This is based on the fact that at term birth, the major sulci and gyri are already present in the neonates [13]. The pattern of the major sulci and gyri are generally preserved but are fine-tuned during brain development [14]. Specifically, the cortical convolutions emerge in the late gestation before birth [15], with extensive folding occurs during the third trimester

<sup>3</sup> <http://neobrain12.isi.uu.nl>

<sup>4</sup> <http://mrbrains13.isi.uu.nl>

<sup>5</sup> <http://www.isles-challenge.org>

<sup>6</sup> <https://www.med.upenn.edu/sbia/brats2017/data.html>

<sup>7</sup> <http://iseg2017.web.unc.edu>

[16, 17]. At term birth, although the brain is only one-third of adult volume [18], the major sulci and gyri present in the adult are already established [13]. Second, based on the obtained initial automatic segmentation, manual editing was performed, under the guidance of an experienced neuroradiologist (Dr. Valerie Jewells, UNC-Chapel Hill), to correct segmentation errors (based on both T1w and T2w MR images) and geometric defects by using ITK-SNAP, with the help of surface rendering. For example, if there is a hole/handle in the surface, we will first localize the related slices, and then check the segmentation maps of both T1w and T2w images to determine whether to fill the hole or cut the handle. Generally, it took almost one week for correcting one subject. Fig. 2 shows an example of a 6-month infant subject with T1w and T2w MR images, and manual labels of WM, GM and CSF, where WM includes unmyelinated white matter and myelinated white matter; GM includes cortical gray matter and subcortical gray matter; and CSF includes the ventricles and cerebrospinal fluid in the extracerebral space. Finally, 10 infant subjects with manual labels are provided for training and 13 infant subjects with manual labels are provided for testing. Note that the manual labels of testing subjects are not provided to the participants for fair comparison. All testing subjects were segmented off site and uploaded for evaluation.

### III. EVALUATION

To evaluate the performance of different methods, we use Dice coefficient (DICE), average surface distance (ASD), and modified Hausdorff distance (HD95) as evaluation metrics to evaluate the accuracy.

#### • DICE

$$\text{DICE} = \frac{2|A \cap B|}{|A| + |B|}$$

where  $A$  and  $B$  denote the binary segmentation labels generated manually and computationally, respectively,  $|A|$  denotes the number of positive elements in the binary segmentation  $A$ , and  $|A \cap B|$  is the number of shared positive elements by  $A$  and  $B$ .

#### • ASD

$$\text{ASD} = \frac{1}{2} \left( \frac{\sum_{V_i \in S_A} \min_{V_j \in S_B} d(V_i, V_j)}{\sum_{V_i \in S_A} 1} + \frac{\sum_{V_j \in S_B} \min_{V_i \in S_A} d(V_j, V_i)}{\sum_{V_j \in S_B} 1} \right)$$

where  $S_A$  is the surface of the ground-truth label map,  $S_B$  is the surface of the automatically segmented label map, and  $d(V_j, V_i)$  indicates the Euclidean distance from vertex  $V_j$  to the vertex  $V_i$ .

#### • HD95

$$\text{HD}(C, D) = \max(h(C, D), h(D, C))$$

where  $C$  and  $D$  are the two sets of vertices identified manually and computationally, respectively, for one tissue class of a subject.  $h(C, D)$  is given by:

$$h(C, D) = \max_{c \in C} \max_{d \in D} \|c - d\|$$

The modified Hausdorff distance is defined as the 95th-percentile Hausdorff distance (HD95).

### IV. METHODS AND IMPLEMENTATIONS

First, we give an overview of all the participants of the iSeg-2017 Challenge, along with a very short description of each participating approach. A total of 21 teams successfully submitted their results to iSeg-2017 before the official deadline. Please refer to Appendix A Table I<sup>8</sup>, in which we describe all the participating teams with affiliations and features used in their methods. In Appendix A Table II<sup>8</sup>, we summarize the performance of all these teams in terms of DICE, ASD and HD95. An interesting finding is that 20 out of 21 teams employed convolutional neural networks for segmentation, while only 1 team utilized a classic atlas-based segmentation method. Among those 20 teams using convolutional neural networks, 8 teams adopted the U-Net architecture [19]. As explained earlier, we will review only the 8 top-ranked methods according to these metrics.

#### A. MSL\_SKKU: Media System Laboratory at Sungkyunkwan University (SKKU), Korea [20]

Bui et al. extended the densely connected convolutional network [21] to deal with segmentation of 6-month infant brain MRI [20]. By concatenating information from shallow to deep dense blocks, the proposed network allows capturing multiple contextual information and yields accurate segmentation results. Their proposed network architecture for infant brain segmentation is shown in Fig. 3.

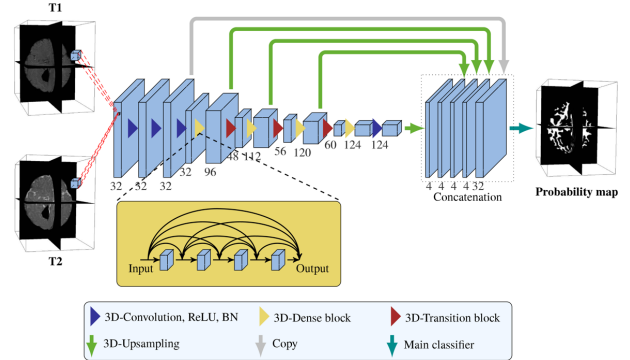


Fig. 3. 3D densely convolutional network architecture for infant brain segmentation.

The network consists of two paths: 1) the down-sampling path and 2) the up-sampling path. The down-sampling path includes four dense blocks. Each dense block comprises of four  $3 \times 3 \times 3$  convolutional kernels, each of which is preceded by a batch normalization (BN) layer [22] and a rectified linear unit (ReLU) nonlinearity [23]. A bottleneck layer is introduced before each  $3 \times 3 \times 3$  convolution to improve computational efficiency. They use a dropout layer [24] with the dropout rate of 0.2 after each  $3 \times 3 \times 3$  convolution layer to avoid over-fitting. Between two contiguous dense blocks, a transition block that has  $1 \times 1 \times 1$  convolution with the compression rate of half and a convolution layer of stride 2 is used to reduce the feature map resolutions while preserving the spatial information. In the up-sampling path, the 3D-upsampling operators are used to recover the input resolution. In particular, the shallower layers provide fine

<sup>8</sup> Supplementary materials are available in Appendix A.



output maps, while the deeper layers contain the coarse output maps [25]. To combine multiple levels of contextual information, up-sampling is performed after each dense block and then those up-sampled feature maps are concatenated. A classifier consisting of a  $1 \times 1 \times 1$  convolution is used to classify the concatenated feature maps into target classes. In total, this network has 47 layers with 1.55 million learnable parameters.

In the implementation, T1w and T2w images were normalized to zero mean and unit variance before inputting them into the network. Due to the limited GPU memory, sub-volume samples of size  $64 \times 64 \times 64$  were used as input of the network. The network was trained with Adam [26] with a mini-batch size of 4. The weights were initialized as in He *et al.* [27]. The learning rate was initially set to 0.0002 and was decreased by a factor of  $\gamma = 0.1$  every 50,000 iterations. Weight decay of 0.0005 and a *momentum* of 0.97 were set up for the network. The final segmentation results were obtained using the majority voting strategy from the predictions of the overlapped sub-volumes with stride of  $8 \times 8 \times 8$ . It took about 2 days for training and 5 minutes for segmenting each subject on a TitanX Pascal GPU and Caffe framework [19, 28].

#### B. LIVIA: Laboratory for Image, Vision and Artificial Intelligence (LIVIA), at the École de technologie supérieure (ETS) in Montreal [29]

Inspired by the recent success of dense networks in image segmentation problems, Dolz et al. proposed an ensemble of semi-dense deep architectures to segment 6-month infant brain MRI [29]. In this novel architecture called SemiDenseNet, the outputs of all convolutional layers are connected directly to the last block of the network. This semi-dense connectivity brings some advantages: 1) efficient propagation of gradients during training, and 2) reducing the number of trainable parameters.

Their proposed method (Fig. 4) extends the recent deep architecture proposed in [30], which is composed of many convolutional layers, each containing several 3D convolution filters. To avoid losing resolution when down-sampling the data, the proposed architecture is a fully convolutional network (FCN) without any pooling operations. In addition, multi-scale context is modeled by embedding outputs from all the layers into a dense feature map that is provided to the first fully connected layer, which gives to the architecture the appearance of a semi-dense CNN. A notable difference of the proposed approach with respect to the most existing works is the adopted sampling strategy. Instead of employing a whole 3D MR scan as the input, they sub-sample the whole image into smaller sub-volumes, which are then fed into the network. This allows: 1) avoiding memory issue if pooling is not employed and 2) removing data augmentation for training, since a high number of samples can be extracted from each image. Further, to achieve a more robust segmentation, an ensemble of several architectures is employed to combine their outputs via a majority voting strategy.

The proposed *SemiDenseNet* is composed of 13 layers in total: 9 convolutional layers in each path, 3 fully-connected layers, and the classification layer. The number of kernels

(with the size of  $3 \times 3 \times 3$ ) in each convolutional layer, from shallow to deeper, is 25, 25, 25, 50, 50, 50, 75, 75 and 75, respectively. The fully-connected layers are composed of 400, 200 and 150 hidden units, respectively, followed by a final classification layer. To preserve spatial resolution, a unit stride is used for all convolutional layers. Each convolutional block is composed by a batch normalization step followed by a Parametric Rectified Linear Unit (PReLU) and several convolutional filters in the convolutional layers. Further, in the fully convolutional connected layers, dropout is employed right after PReLU activations. The optimization of network parameters is performed via RMSprop optimizer. *Momentum* was set to 0.6 and the initial learning rate to 0.001, reduced by a factor of 2 after every 5 epochs (starting from epoch 10). Weights in layer  $l$  were initialized based on a zero-mean Gaussian distribution of standard deviation  $2/nl$ , where  $nl$  denotes the number of connections to units in that layer. The proposed 3D FCN was trained for 30 epochs, each composed of 20 subepochs. At each subepoch, a total of 1000 samples were randomly selected from the training images, and processed in batches of size 20. An ensemble composed by 10 identical CNNs was employed, each trained with a different combination of subjects. No data augmentation was employed to increase the size of the dataset. Experiments were performed in a computational server equipped with a NVIDIA Tesla P100 GPU with 16 GB of RAM memory. Training the proposed network took around 25 min per epoch, and around 13 hours to have a single CNN. Segmentation of a whole 3D MR scan was performed in 10 seconds per CNN model in average.

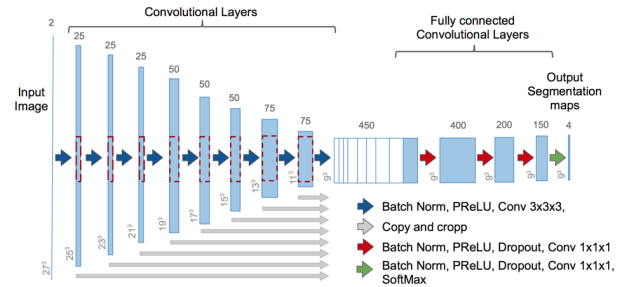


Fig. 4. Architecture of the proposed SemiDenseNet, which takes as the input sub-patches of size  $27 \times 27 \times 27$  from T1w and T2w images and provides segmentation maps of size  $9 \times 9 \times 9$ .

#### C. Bern\_IPMI: Information Processing in Medical Intervention Lab., University of Bern, Switzerland [31]

Zeng and Zheng proposed a two-stage, 3D fully convolutional networks (3DFCN)-based method for segmentation of 6-month infant brain MRI [31]. In order to alleviate the potential gradient vanishing problem during training, they designed multi-scale deep supervision. Moreover, context information was used to further improve the performance.

Fig. 5 illustrates their proposed two-stage method. Both 3DFCN-1 and 3DFCN-2 adopt an encoder (contracting path)-decoder (expansive path) structure [32]. More specifically, 3DFCN-1 is used in the first stage to learn the probability map of each brain tissue from multimodal MR images (T1w and T2w images). An initial segmentation of



different brain tissues is then obtained from the probability map, which further allows us to compute a distance map for each tissue [33]. The computed distance maps can be used to model the spatial context information. At the second stage, 3DFCN-2 is employed to get the final segmentation by using both the spatial context information and the multimodal MR images. To effectively integrate multimodal information, separate encoder paths are constructed for different modalities and then their outputs of the encoder paths are concatenated at the beginning of the expansive path such that the decoder can fuse complementary information from different sources. At both stages, long and short skip connections are employed to recover spatial context lose in the contracting encoder. See Fig. 5 for details. For 3DFCN-2, two down-scaled branch classifiers are further injected into the networks in addition to the classifier of the main network. By doing this, segmentation is performed at multiple output layers. As a result, classifiers in different scales can take advantage of multi-scale context.

Their proposed method was implemented in Python using TensorFlow framework and trained on a desktop with a 3.6 GHz Intel® i7 CPU and a GTX 1080 Ti graphics card with 11 GB GPU memory. In order to enlarge the training samples, data augmentation was utilized. Specifically, each training data was rotated for (90, 180, 270) degrees around the y-axis of the image and flipped horizontally (by taking the z axis as the vertical direction). The network was trained for 10,000 iterations. All weights were updated by the stochastic gradient descent algorithm (*momentum*=0.9, *weight\_decay*=0.005). Learning rate was initialized as  $1 \times 10^{-3}$  and halved by every 3,000 times. After training, the proposed method took about 8 seconds in average to segment one subject.

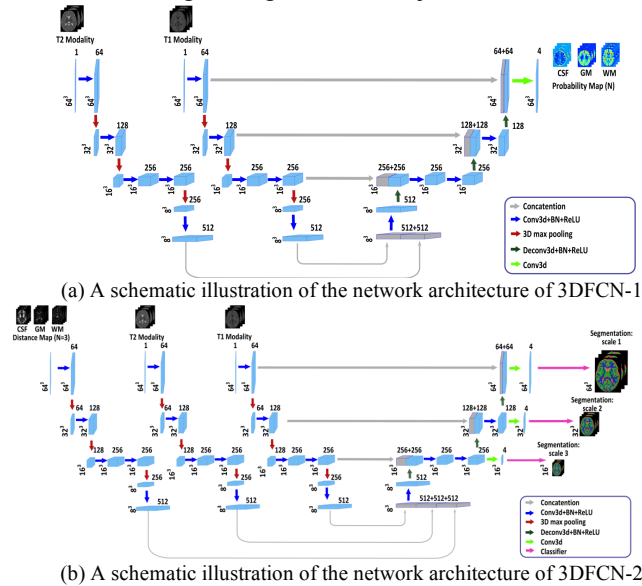


Fig. 5. A schematic illustration of the proposed two-stage method, consisting of (a) 3DFCN-1 at stage one and (b) 3DFCN-2 at stage two. For each block, the numbers above represent the number of feature stacks, and the numbers at the left side indicate data size.

*D. TU/e IMAG/e: author list; Medical Image Analysis Group (IMAG/e) of Eindhoven University of Technology (TU/e) [34]*

A convolutional neuronal network was used for the segmentation of 6-month infant brain MRI into WM, GM and CSF [34]. Unlike previous work [35], the network does not include pooling layers, but uses dilated convolutions to achieve a large receptive field using a limited number of trainable weights.

The method combines 2D triplanar and 3D input using four network branches (Fig. 6). All network branches use the T1w and T2w images as 2-channel input. The triplanar input is included in three branches with dilated 2D convolutions. Each of these branches consists of 7 layers of  $3 \times 3$  convolutions with increasing dilation factors, resulting in a receptive field of  $67 \times 67$  [36], as previously also used for cardiac segmentation [37] and adult brain MRI segmentation [34]. The 3D input is included in the fourth branch that consists of 12 layers of  $3 \times 3 \times 3$  convolutions, resulting in a receptive field of  $25 \times 25 \times 25$ . The output features from the four branches are concatenated and combined in the output layer with  $1 \times 1$  convolutions.

Batch normalization and ReLUs were used throughout. Dropout was used before the output layer. The network was trained with Adam based on the cross-entropy loss, using mini-batches of 200 or 300 samples in 10 epochs of 50,000 random samples per class per training image. The network was trained with a patch-based approach, randomly sampling from all images in the training set. During the testing, arbitrarily sized inputs can be used, because of the fully convolutional nature of all four branches. The method took about 1 minute to segment a full image on a NVIDIA Titan X Pascal GPU. The segmentation results were obtained without any data augmentation. Data augmentation could possibly improve the results in scenarios not well represented in the training set.

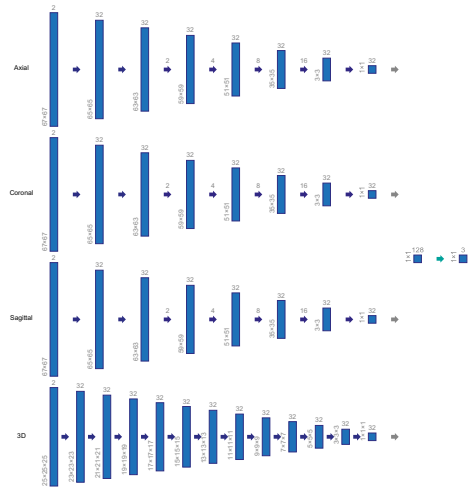


Fig. 6. Network architecture. The colors of the arrows indicate, from left to right:  $3 \times 3$  or  $3 \times 3 \times 3$  convolutions, concatenation, and  $1 \times 1$  convolutions. Dilation factors are shown above the arrows. During the training, single voxels are used as output. During the testing, arbitrarily sized outputs can be used, because of the fully convolutional nature of the network.

*E. UPF\_Simbiosys: Simbiosys research lab at Universitat Pompeu Fabra (UPF), Barcelona [38]*

There exist many segmentation approaches, such as multi-atlas label fusion [39] and learning-based methods [40, 41]. Each method has its own strength, and different segmentation approaches may potentially complement each other. The motivation for the proposed method is to combine the strengths of complementary methods in a cascaded fashion.

The pipeline of the method is shown in Fig. 7. The 0-level of the cascade segments the multi-modal (T1w and T2w) input images independently with joint label fusion (JLF) [39]. The estimated probability maps in level-0, along with the original images, are inputted to the level-1 of the cascade. In level-1, first, multi-scale features are extracted from both input images and probability maps of level-0. Image features consist of 1) Gaussian, 2) Laplacian-of-Gaussian, and 3) gradient magnitude images convolved with Gaussians at multiple scales for each modality. Probability features are obtained by convolving the level-0 probability maps with Gaussians at multiple scales. The multi-scale image and probability-map features are fed into a SVM classifier for outputting the final estimated label map. Each sample of the SVM classifier is composed of the features extracted from each voxel. The SVM classifier is trained during the training phase using the features extracted from the training set.

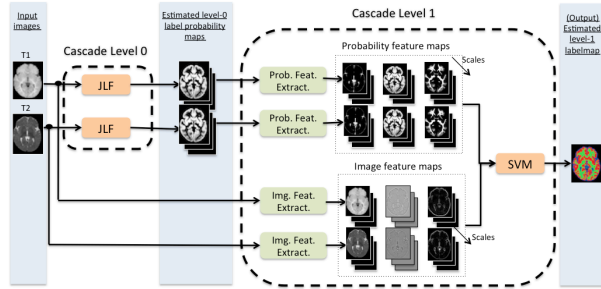


Fig. 7. Dashed blocks correspond to the different levels of the cascade. Blue columns denote input, intermediate output, and final results. Rounded rectangles denote segmentation methods (orange) and feature extraction processes (green), respectively.

Pre-processing steps include 1) histogram matching of all the images to the UNC 1-year-old infant template [11], and 2) non-rigid registration to the same template using ANTs [42]. Pair-wise registrations for multi-atlas JLF are computed by concatenating registrations through the template. No post-processing steps are applied. The parameters for the segmentation methods in each level (i.e., JLF and SVM) are chosen by cross-validation in the training set. Specifically, for JLF, the patch radius is set to be 2 for both modalities and the search window is set to be 7 and 5 for T1w and T2w images, respectively. For SVM, we set the regularization constant to  $C=5$ , use an RBF kernel, and normalize the features to zero-mean and unit standard deviation. The computational time for segmenting each subject is  $\sim 30$  minutes.

The performance of the SVM classifier in level-1 is highly influenced by the features derived from JLF in level-0. This suggests the advantage of combining multiple complementary methods in the proposed cascaded scheme. A slight drop in

performance is experienced by adding an extra layer in the cascade by the level-1 outputs using as the input, so the two-levels scheme is kept as the final model. Among different combination strategies, the proposed cascaded scheme performed better than an alternative ensembling strategy [43].

*F. NeuroMTL: Montreal Neurological Institute, McGill University, Montreal QC Canada<sup>9</sup>[44]*

First, an extended training dataset was created by applying existing tissue classification to scans from the longitudinal dataset of infants at-risk of autism and control subject in the Infant Brain Imaging Study (IBIS) [45] where scans of 24-month old infants for whom 6 and 12 month scans were available and had T1w and T2w scans acquired at all time points ( $n=216$ ).

Tissue classification method is shown in Fig. 8: i) An unbiased population average of T1w scans for each age group (6 months, 12 months and 24 months) was created [46]. ii) The group average for the 24-month old scans was manually segmented into areas of high probability of WM, GM and CSF. iii) All 24-month-old T1w scans were non-linearly registered to the template, and tissue priors from the template were transformed to the space of each subject's scan. iv) An expectation-maximization algorithm was run to obtain tissue classification. v) Longitudinal non-linear registrations between scans at 6 and 12 months and then between 12 and 24 months were performed using ANTs with mutual information [42], using both T1w and T2w scans. Using these registration transformations, tissue classification maps from 24 months were transformed to the 6-month scans. Segmentations from the 24-month scans were propagated back to the 6-month scans via non-linear registration. Then, a 3D U-Net [19] was trained in two stages with the extra dataset to automatically segment healthy tissues. U-Net with 5 downsampling and upsampling blocks with skip connections was trained on  $80^3$  image patches for tissue classification, with the parameters listed in Table I. Each block contained two convolutional layers with ReLU activations, with  $5 \times 5 \times 5$  convolution layers in the first two blocks,  $3 \times 3 \times 3$  convolution layers in the next two blocks, and a combination of  $3 \times 3 \times 3$  and  $1 \times 1 \times 1$  convolution layers in the fifth block, with max pooling at each block. Additionally, a  $3 \times 3 \times 3$  convolution layer with 64 input and output channels was added, followed by a  $1 \times 1 \times 1$  convolutional layer with 64 input and 32 output channels and then another  $1 \times 1 \times 1$  layer with 32 input and 4 output channels with dropout, optimizing categorical cross-entropy with Adam. The output patch was cropped to  $64^3$  to remove edge effects. Training was done in two stages, first on the IBIS dataset, and then fine-tuned on the iSeg-2017 challenge data ( $n=10$ ).

<sup>9</sup> Fonov et al. acknowledged imaging data was collected as part of the Infant Brain Imaging Study (IBIS). Fonov et al. also thank IBIS children and families for their ongoing participation in this longitudinal study.

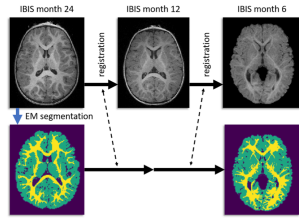


Fig. 8. Automatic segmentation of 6-month old infant MRI data.

All experiments were performed on a computer with Xeon CPU E5-2620 v4 @ 2.10GHz with 64GB of ram and NVIDIA Titan-X GPU, with deep-net implemented in Torch7. Training on ACE-IBIS dataset took approximately 32 hours (10000 mini-batches), and final training on iSeg-2017 data took 11 hours (4000 mini-batches). Application on a single subject, using GPU, took 8 seconds.

TABLE I. PARAMETERS OF 3D U-NET.

Layer	Input Channels	Output Channels	Convolution kernel 2	Convolution kernel 2	Upsampling kernel
1	4	64	5x5x5	5x5x5	5x5x5
2	16	64	5x5x5	5x5x5	3x3x3
3	16	64	3x3x3	3x3x3	3x3x3
4	16	64	3x3x3	3x3x3	3x3x3
5	32	64	1x1x1	3x3x3	-

#### G. UPC\_DLMI: Signal Theory and Communications Department, Universitat Politècnica de Catalunya, Barcelona

Casamitjana et al. proposed a convolutional neural network, named Augmented V-Net (Fig. 9), which is an extension of the V-Net architecture [47]. The main changes with respect to the original V-Net model can be summarized as follows:

- *Augmented path*: An upsampled version of the input is used to exploit high resolution features. This is done by upsampling by repetition the input (factor of 2) and stacking several convolutional layers after the upsampling. The resultant features are concatenated in the last layers.
- *Modified residual connections*: The residual connections are reformulated such that the propagation of the input signal through the network is minimally modified.
- *Mask*: A mask is used before the final prediction in order to constrain the network to train on relevant voxels.
- *Input concatenation*: The raw input image is used as feature map in the last stages of the network.

The key part of the network is the augmented path, which has been shown to boost the performance of the standard V-Net for the infant brain segmentation task. It provides high-resolution features by keeping small filter sizes and adding redundancy in the input, helping to detect finer regions such as boundaries. Later in the network, the authors use the input image as raw features, since voxel's intensities already contain valuable information. Finally, the mask is used to train/predict only on voxels of brain tissue.

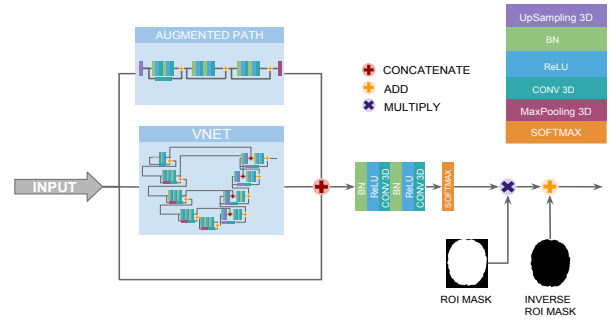


Fig. 9. Augmented V-Net. It builds upon the concatenation of the V-Net core network [47] with an augmented path with higher resolution. Augmented V-Net uses a ROI-mask to train only in brain tissue voxels. Layer types are color-coded as shown in the top-right corner.

T1w and T2w MRIs are used as input images. Both are normalized to zero mean and unit variance. From the normalized T1w image, a mask is created to mask out background voxels. When training such a big and deep network, there are two main problems: GPU memory constraints and the scarcity of data. Patch-wise training arises as a possible solution for the first issue. The memory required to train Augmented V-Net does not allow using dense-training, which is also discouraged when data is scarce. Then, patch-wise training is the only solution. Larger patches are preferred because they can encode localization features (brain structures) across the network, while smaller patches allow increasing the batch size in the optimization process. The authors finally choose patches of size  $64 \times 64 \times 64$  and sample uniformly across the brain, forcing the central voxel to belong to brain tissue (WM, GM and CSF).

Interestingly, faster and better convergence was obtained with lower generalization error, since batch sizes were composed of more subjects, relaxing the problem of data scarcity. The authors used data augmentation to increase the size of the training set, by making sagittal reflections of each subject. Other reflections have been shown to produce worse results, and no other datasets were used to train the network. In the optimization process, they used Adam optimizer with initial learning rate of  $lr=0.0005$ . The loss function used was the weighted cross-entropy, where loss weights were computed as the normalized inverse of the class frequency. At inference time, the whole subject can be used as input for the trained model, performing dense inference and using the mask to indicate brain tissue voxels. The method is fully automatic, taking from 5 to 7 seconds to process one subject.

#### H. LRDE: EPITA Research & Development Laboratory [48]

Xu et al.'s method is an extension from single modality to multi-modality of the authors' previous work on neonatal infant brain MRI segmentation [48]. This automatic method uses fully convolutional network (FCN) and transfer learning (see details in Fig. 10), and is very fast: the segmentation of a whole volume only takes a few seconds. The core part of the 16 layers VGG network [43] is used. This very efficient network, pre-trained on millions of 2D color natural images in ImageNet (for image classification purpose), and fine-tuned with the MRI training dataset. The key contribution is to show



how to build 2D color images from a 3D MRI volume, so that VGG effectively gives state-of-the-art segmentation results.

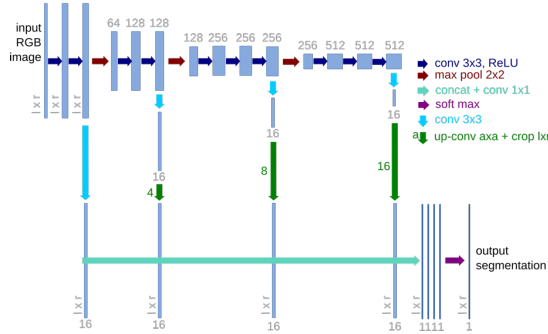


Fig. 10. Visualization of the proposed segmentation network

The combination of the T1w and T2w slices to obtain a set of 2D color (RGB) images is very simple. For each slice (indexed by  $n$ ), the fake color image is constructed in such a way that the “green” channel is the T2w slice  $n$ , and the red and the blue are T1w slices respectively at indices  $n-1$  and  $n+1$ . Each 2D color image thus forms a *3D-like* representation of a part (3 consecutive slices) of the MR volume. This representation enables incorporating some 3D information, while avoiding the expensive computational and memory requirements of fully 3D CNN. For this specific application, the fully connected layers at the end of VGG network are discarded; only the 4 stages of convolutional parts called “base network” are retained. This base network is composed of convolutional layers, ReLU layers and max pooling layers between two successive stages. The three max pooling layers divide the base network into four stages of fine to coarse feature maps. A stack of specialized layers is obtained, 1 from each stage, and a softmax function yields the segmentation result.

TABLE II. SOURCE CODES FROM TOP TEAMS IN ISEG-2017.

TEAM	METHOD	LINK
<i>MSL_SKKU</i>	3D, DenseNet	<a href="https://github.com/tbuikr/3D_DenseSeg">https://github.com/tbuikr/3D_DenseSeg</a>
<i>LIVIA</i>	3D, CNN+FC	<a href="https://github.com/josedolz/SemiDenseNet">https://github.com/josedolz/SemiDenseNet</a>
<i>Bern_IPMI</i>	3D, Two stages, double-armed U-Net	<a href="https://github.com/zengguodong/iSeg_Bern_IPMI">https://github.com/zengguodong/iSeg_Bern_IPMI</a>
<i>TU/e IMAG/e</i>	3D, CNN with dilated convolutions	<a href="https://github.com/pimmoeskops/iSeg_dilatedCNN">https://github.com/pimmoeskops/iSeg_dilatedCNN</a>
<i>NeuroMTL</i>	3D, U-Net	<a href="https://github.com/vfonov/NeuroMTL_iSEG">https://github.com/vfonov/NeuroMTL_iSEG</a>
<i>UPC_DLMI</i>	3D, Augmented U-Net	<a href="https://github.com/imatge-upc/segmentation_DLMI/">https://github.com/imatge-upc/segmentation_DLMI/</a>
<i>LRDE</i>	2D, Pretrained VGG16+FCN	<a href="https://www.lrde.epita.fr/wiki/NeoBrainSeg">https://www.lrde.epita.fr/wiki/NeoBrainSeg</a>

## V. DISCUSSION

Based on Section IV, most of the well-performed teams (7 teams out of 8) adopted deep learning based algorithms. Moreover, most of the deep learning related algorithms are based on 3D U-Net (or U-Net-like structures). Thanks to the use of GPUs, most of these algorithms have inference times between 5-10 seconds for a whole MRI scan. The only non-deep learning based method is developed by Sanroma et al. (*UPF\_simbiosys*), which employs a multi-atlas based method followed by an SVM to design a cascade learning segmentation algorithm.

Before creating the set of 2D color images, a pre-processing of the T1w and T2w sequences was performed, which consists of: 1) shifting the voxel values of the MRI volumes to center their histograms on their maximal histogram value, and 2) quantizing the voxel values on 8bit (values lower than 0 and greater than 255 are saturated). For the training, the classical data augmentation strategy by scaling and rotating images were adopted. 2D images were then computed for each volume of the augmented training base using the pre-processed T1w and T2w slices as described before. The network was fine-tuned for the first 50K iterations with a learning rate of  $lr=10^{-8}$ , and the last 100K with a smaller learning rate ( $lr=10^{-10}$ ). Stochastic gradient descent was employed to minimize the loss function with *momentum*=0.99 for the first 50K iterations and 0.999 for the next 100k, and *weight\_decay*=0.0005. The loss function was averaged over 20 images. During test, the runtime on a 3D volume was 1.8 seconds on average; note that this included the pre-processing step, the computation of the set of 2D color input images, and after inference, the reconstruction of a 3D volume (the expected segmentation output) by stacking the set of 2D output images.

### I. Source Codes

A proactive goal of this paper is to encourage authors to make their codes publicly available for reproducible research. By far, most of teams have shared their codes, as summarized in Table. II. For readers who seek to come up to speed with deep learning, these codes can be also served as good starting points to understand how deep learning algorithms can be implemented for image segmentation.

### V.I. Evaluation on the whole brain

We first evaluate the performance in terms of the whole brain. Fig. 11 reports performances of 8 teams using DICE, HD95 and ASD, with box-plots. Besides medians, means are also indicated by the black diamonds. To see if any method produces a significantly better result, we calculated Wilcoxon two tailed-test, as shown in Table. III with all-against-all diagram in terms of metrics (DICE, HD95 and ASD). Interesting, we did not find any method could achieve a strong statistically significant difference ( $p$ -value<0.01) with other methods, in terms of WM, GM and CSF using any metric (DICE, HD95 and ASD). For example, in terms of WM DICE, we found *MSL\_SKKU* has the highest median, but there are no

strong statistically significant difference with *LIVIA*, and also with *Bern\_IPMI*. In terms of WM HD95, *LRDE* has the lowest median, but there is no any statistically significant difference with any other methods. For WM ASD, *MSL\_SKKU* has the lowest median, but there is no statistically significant difference with *LIVIA* and *Bern\_IPMI*. For GM DICE, we found *MSL\_SKKU* has the highest median, but there is no

statistically significant difference with *MSL\_SKKU* and *Bern\_IPMI*. For GM HD95, *MSL\_SKKU* has the lowest median, but there is no strong statistical significance with any other methods. For GM ASD, *MSK\_SKKU* has the lowest median, but there is no statistically significant difference with *LIVIA*.

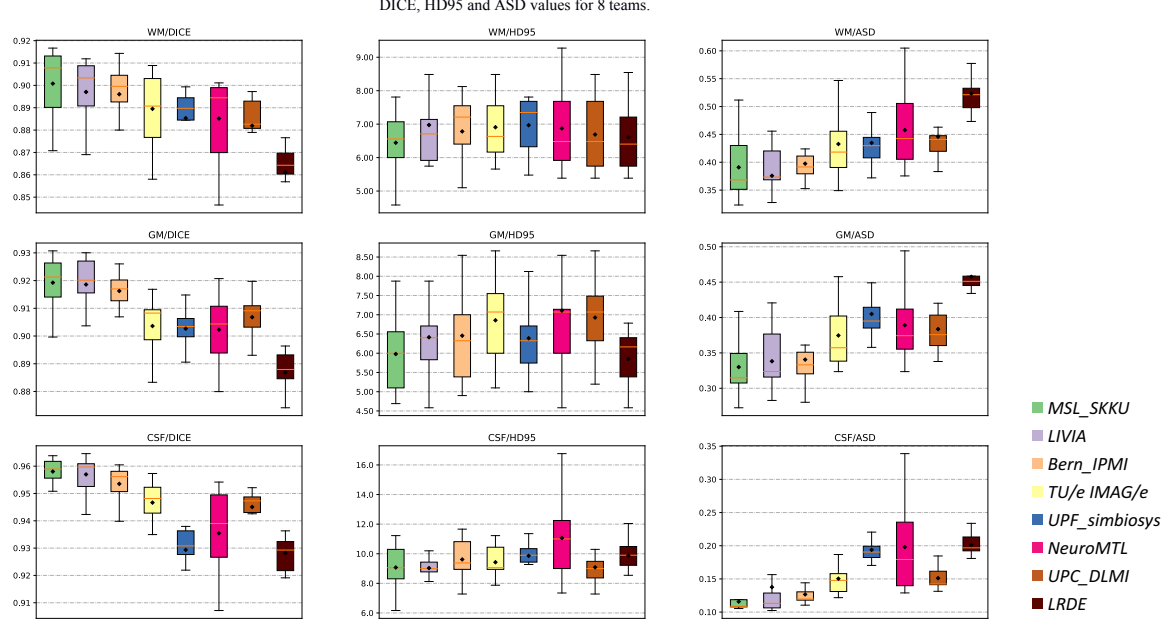


Fig. 11. Performances of 8 teams in terms of DICE, HD95 and ASD, with box-plots. Besides medians, means are also indicated by the dark dots.

TABLE III. P-VALUES BY PERFORMING WILCOXON SIGNED-RANK TEST

<sup>+</sup> Denotes weak statistical significance ( $p$ -value  $< 0.05$ ).

<sup>++</sup> Denotes strong statistical significance ( $p$ -value  $< 0.01$ ).

WM/DICE								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.033 <sup>+</sup>	0.023 <sup>+</sup>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>
<i>LIVIA</i>	0.033 <sup>+</sup>	N/A	0.507	0.005	0.001 <sup>+</sup>	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>
<i>Bern_IPMI</i>	0.023 <sup>+</sup>	0.507	N/A	0.028	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.001 <sup>++</sup>
<i>TU/e IMAG/e</i>	0.001 <sup>++</sup>	0.005 <sup>++</sup>	0.028 <sup>+</sup>	N/A	0.221	0.075	0.023 <sup>+</sup>	0.001 <sup>++</sup>
<i>UPF_simbiosys</i>	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.221	N/A	0.917	0.033 <sup>+</sup>	0.001 <sup>++</sup>
<i>NeuroMTL</i>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.075	0.917	N/A	0.249	0.001 <sup>++</sup>
<i>UPC_DLMI</i>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.023 <sup>+</sup>	0.033 <sup>+</sup>	0.249	N/A	0.001 <sup>++</sup>
<i>LRDE</i>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	N/A

WM/HD95 (mm)								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.347	0.182	0.311	0.019 <sup>+</sup>	0.196	0.534	0.753
<i>LIVIA</i>	0.347	0.000	0.507	0.530	0.972	0.480	0.388	0.117
<i>Bern_IPMI</i>	0.182	0.507	N/A	0.807	0.480	0.638	0.753	0.695
<i>TU/e IMAG/e</i>	0.311	0.530	0.807	N/A	0.917	0.917	0.345	0.158
<i>UPF_simbiosys</i>	0.019 <sup>+</sup>	0.972	0.480	0.917	N/A	0.530	0.374	0.221
<i>NeuroMTL</i>	0.196	0.480	0.638	0.917	0.530	N/A	0.221	0.131
<i>UPC_DLMI</i>	0.534	0.388	0.753	0.345	0.374	0.221	N/A	0.916

<i>LRDE</i>	0.753	0.117	0.695	0.158	0.221	0.131	0.916	N/A
-------------	-------	-------	-------	-------	-------	-------	-------	-----

WM/ASD (mm)								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.173	0.311	0.001 <sup>††</sup>	0.006 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>
<i>LIVIA</i>	0.173	N/A	0.917	0.003 <sup>††</sup>	0.003 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>
<i>Bern_IPMI</i>	0.311	0.917	N/A	0.004 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>
<i>TU/e IMAG/e</i>	0.001 <sup>††</sup>	0.003 <sup>††</sup>	0.004 <sup>††</sup>	N/A	0.807	0.016 <sup>†</sup>	0.196	0.001 <sup>††</sup>
<i>UPF_simbiosys</i>	0.006 <sup>††</sup>	0.003 <sup>††</sup>	0.002 <sup>††</sup>	0.807	N/A	0.196	0.023 <sup>†</sup>	0.001 <sup>††</sup>
<i>NeuroMTL</i>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.016 <sup>†</sup>	0.196	N/A	0.221	0.001 <sup>††</sup>
<i>UPC_DLMI</i>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.196	0.023 <sup>†</sup>	0.221	N/A	0.001 <sup>††</sup>
<i>LRDE</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	N/A

GM/DICE								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.650	0.055	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>
<i>LIVIA</i>	0.650	N/A	0.087	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>
<i>Bern_IPMI</i>	0.055	0.087	N/A	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>
<i>TU/e IMAG/e</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	N/A	0.311	0.422	0.249	0.002 <sup>††</sup>
<i>UPF_simbiosys</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.311	N/A	0.861	0.005 <sup>††</sup>	0.001 <sup>††</sup>
<i>NeuroMTL</i>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.422	0.861	N/A	0.133	0.002 <sup>††</sup>
<i>UPC_DLMI</i>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.249	0.005 <sup>††</sup>	0.133	N/A	0.001 <sup>††</sup>
<i>LRDE</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	N/A

GM/HD95 (mm)								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.600	0.169	0.033 <sup>†</sup>	0.093	0.101	0.028 <sup>†</sup>	0.442
<i>LIVIA</i>	0.600	N/A	0.421	0.158	0.937	0.071	0.173	0.074
<i>Bern_IPMI</i>	0.169	0.421	N/A	0.382	0.814	0.650	0.182	0.081
<i>TU/e IMAG/e</i>	0.033 <sup>†</sup>	0.158	0.382	N/A	0.173	0.972	0.875	0.016 <sup>†</sup>
<i>UPF_simbiosys</i>	0.093	0.937	0.814	0.173	N/A	0.196	0.221	0.060
<i>NeuroMTL</i>	0.101	0.071	0.650	0.972	0.196	N/A	1.000	0.008 <sup>††</sup>
<i>UPC_DLMI</i>	0.028 <sup>†</sup>	0.173	0.182	0.875	0.221	1.000	N/A	0.008 <sup>††</sup>
<i>LRDE</i>	0.442	0.074	0.081	0.016b	0.060	0.008 <sup>††</sup>	0.008 <sup>††</sup>	N/A

GM/ASD (mm)								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.116	0.028 <sup>†</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>
<i>LIVIA</i>	0.116	N/A	0.463	0.004 <sup>††</sup>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>
<i>Bern_IPMI</i>	0.028 <sup>†</sup>	0.463	N/A	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>
<i>TU/e IMAG/e</i>	0.001 <sup>††</sup>	0.004 <sup>††</sup>	0.001 <sup>††</sup>	N/A	0.013 <sup>†</sup>	0.075	0.173	0.001 <sup>††</sup>
<i>UPF_simbiosys</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.013 <sup>†</sup>	N/A	0.064	0.002 <sup>††</sup>	0.001 <sup>††</sup>
<i>NeuroMTL</i>	0.001 <sup>††</sup>	0.002 <sup>††</sup>	0.001 <sup>††</sup>	0.075	0.064	N/A	0.507	0.001 <sup>††</sup>
<i>UPC_DLMI</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.173	0.002 <sup>††</sup>	0.507	N/A	0.001 <sup>††</sup>
<i>LRDE</i>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	0.001 <sup>††</sup>	N/A

CSF/DICE								
----------	--	--	--	--	--	--	--	--



TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.552	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>
<i>LIVIA</i>	0.552	N/A	0.009 <sup>++</sup>	0.007 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>
<i>Bern_IPMI</i>	0.002 <sup>++</sup>	0.009 <sup>++</sup>	N/A	0.039 <sup>+</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.001 <sup>++</sup>
<i>TU/e IMAG/e</i>	0.001 <sup>++</sup>	0.007 <sup>++</sup>	0.039 <sup>+</sup>	N/A	0.001 <sup>++</sup>	0.033 <sup>+</sup>	0.650	0.001 <sup>++</sup>
<i>UPF_simbiosys</i>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	N/A	0.133	0.001 <sup>++</sup>	0.116
<i>NeuroMTL</i>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.033 <sup>+</sup>	0.133	N/A	0.075	0.116
<i>UPC_DLMI</i>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.002 <sup>++</sup>	0.650	0.001 <sup>++</sup>	0.075	N/A	0.001 <sup>++</sup>
<i>LRDE</i>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.116	0.116	0.001 <sup>++</sup>	N/A

CSF/HD95 (mm)								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.972	0.169	0.477	0.221	0.041 <sup>+</sup>	0.929	0.133
<i>LIVIA</i>	0.972	N/A	0.064	0.117	0.050	0.011 <sup>+</sup>	0.953	0.012 <sup>+</sup>
<i>Bern_IPMI</i>	0.169	0.064	N/A	0.401	0.196	0.050	0.041 <sup>+</sup>	0.311
<i>TU/e IMAG/e</i>	0.477	0.117	0.401	N/A	0.099	0.021 <sup>+</sup>	0.249	0.173
<i>UPF_simbiosys</i>	0.221	0.050	0.196	0.099	N/A	0.311	0.050	0.754
<i>NeuroMTL</i>	0.041 <sup>+</sup>	0.011 <sup>+</sup>	0.050	0.021 <sup>+</sup>	0.311	N/A	0.010 <sup>+</sup>	0.209
<i>UPC_DLMI</i>	0.929	0.953	0.041b	0.249	0.050	0.010 <sup>+</sup>	N/A	0.006 <sup>++</sup>
<i>LRDE</i>	0.133	0.012 <sup>+</sup>	0.311	0.173	0.754	0.209	0.006 <sup>++</sup>	N/A

CSF/ASD (mm)								
TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.861	0.002 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>
<i>LIVIA</i>	0.861	N/A	0.101	0.055	0.023 <sup>+</sup>	0.019 <sup>+</sup>	0.023 <sup>+</sup>	0.023 <sup>+</sup>
<i>Bern_IPMI</i>	0.002 <sup>++</sup>	0.101	N/A	0.011 <sup>+</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>
<i>TU/e IMAG/e</i>	0.001 <sup>++</sup>	0.055	0.011 <sup>+</sup>	N/A	0.001 <sup>++</sup>	0.039 <sup>+</sup>	0.917	0.001 <sup>++</sup>
<i>UPF_simbiosys</i>	0.001 <sup>++</sup>	0.023 <sup>+</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	N/A	0.701	0.001 <sup>++</sup>	0.033 <sup>+</sup>
<i>NeuroMTL</i>	0.001 <sup>++</sup>	0.019 <sup>+</sup>	0.001 <sup>++</sup>	0.039 <sup>+</sup>	0.701	N/A	0.046 <sup>+</sup>	0.650
<i>UPC_DLMI</i>	0.001 <sup>++</sup>	0.023 <sup>+</sup>	0.001 <sup>++</sup>	0.917	0.001 <sup>++</sup>	0.046 <sup>+</sup>	N/A	0.001 <sup>++</sup>
<i>LRDE</i>	0.001 <sup>++</sup>	0.023 <sup>+</sup>	0.001 <sup>++</sup>	0.001 <sup>++</sup>	0.033 <sup>+</sup>	0.650	0.001 <sup>++</sup>	N/A

## V.II. Evaluation based on 80 ROIs

Besides evaluation in terms of the whole brain, we further evaluate the performances based on 80 ROIs. Specifically, a total of 33 two-year-old subjects were employed as individual atlases (www.brain-development.org) [49]. Each atlas consists of a T1w MR image and a label image of 80 ROIs (excluding cerebellum and brainstem). We first employ FreeSurfer [50] to segment each T1w MR image into WM, GM, and CSF. Then, we register all atlases into each testing subject space based on their segmentations using ANTs [51]. Finally, we employ majority voting to parcellate each testing subject into 80 ROIs. For each ROI, we employed DICE to measure the performance between automatic segmentations and manual segmentation. Average ROI-based DICEs for 8 teams are shown in Table. IV. The large number of ROIs reduced the organizers' willingness to report  $p$ -values for each ROI. However, to better interpret these ROI-based evaluations, we have generated error maps for each method, as shown in Fig. 12. They were estimated by aligning all the error maps from 13 testing subjects to a 6-month template [52]. The higher

value of error map, the higher probability for miss-classification. From all these error maps, we can see all methods consistently produce small errors in the subcortical regions while large errors in the cortical regions, which is actually consistent with the fact that tissue contrast is much lower in the cortical regions than subcortical regions. Average error map for all 8 methods were further generated, as shown in the right bottom of Fig. 12. The most error-prone ROI regions are straight gyrus, lingual gyrus, and medial orbital gyrus. These regions are also consistently confirmed with Table V, where DICEs of these ROIs were relatively low, with around 0.84. By contrast, the Dice ratios of subcortical regions, such as putamen and thalamus, are as high as 0.94.

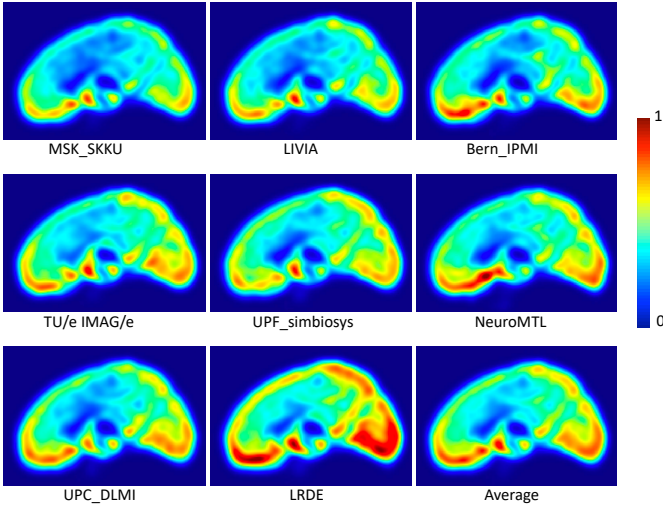


Fig. 12. Error maps: All methods produce small errors in the subcortical regions while large errors in the cortical regions. The most error-prone regions are straight gyrus, lingual gyrus, and medial orbital gyrus. Average error map for all 8 teams is shown in the right bottom.

### V.III. Evaluation based on gyral landmark curves

To better reflect the accuracy of the methods on the gyral crests, we further measure the distance of gyral landmark curves on the cortical surfaces. Large curve distance error indicates that the gyral crest is poorly resolved. We selected two major gyri, i.e., the superior temporal gyral curve and the postcentral gyral curve, as the landmarks to measure the accuracy. We manually labeled two sets of gyral curves on the inner cortical surfaces from different tissue segmentation results. One typical example is shown as in Fig. 13, in which the curves were delineated by the experts on the superior temporal gyrus and postcentral gyrus, the white curve indicate the ground truth, and the colored lines indicates different method respectively. We employed HD95 to calculate the

curves distance, with median HD95 over 13 testing subjects shown in Fig. 14.  $P$ -values were calculated based on Wilcoxon two-tailed test, as shown in Table. III. We find that *Bern\_IPMI* achieves the lowest median HD95, but there is no statistically significant difference with *MSK\_SKKU*, *LIVIA*, *UPF\_simbiosys*, and *NeuroMTL*.

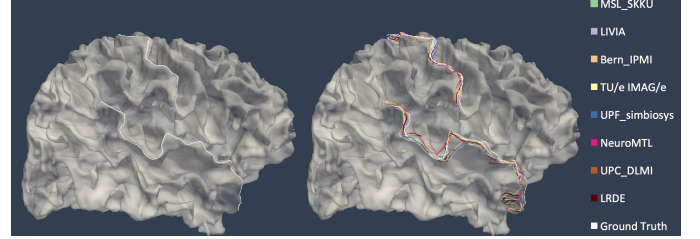


Fig. 13. Evaluations on gyri/sulci for 8 teams. The left one shows the manually labeled postcentral and superior temporal gyral landmark curves of ground truth, while the right one shows the curves of the segment results by 8 different methods compared with ground truth.

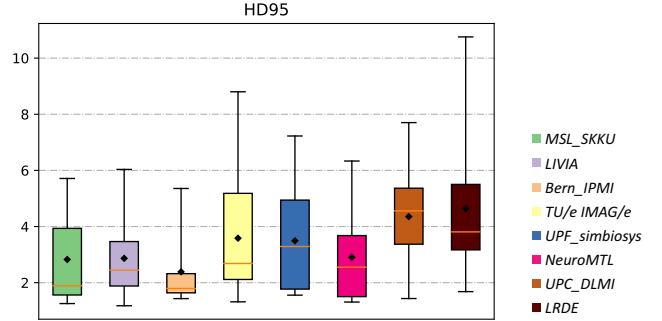


Fig. 14. The boxplot shows HD95 evaluations of 8 different methods on the superior temporal gyrus and the postcentral gyrus of 13 subjects. Besides medians, means are also indicated by the dark dots.

Table. III.  $p$ -values for the two-sided Wilcoxon paired signed-rank test.

<sup>†</sup> Denotes weak statistical significance ( $p$ -value  $< 0.05$ ).

<sup>††</sup> Denotes strong statistical significance ( $p$ -value  $< 0.01$ ).

TEAM	<i>MSL_SKKU</i>	<i>LIVIA</i>	<i>Bern_IPMI</i>	<i>TU/e IMAG/e</i>	<i>UPF_simbiosys</i>	<i>NeuroMTL</i>	<i>UPC_DLMI</i>	<i>LRDE</i>
<i>MSL_SKKU</i>	N/A	0.6221	0.3804	0.2036	0.4697	0.8501	0.0068 <sup>††</sup>	0.0640
<i>LIVIA</i>	0.6221	N/A	0.4238	0.0269	0.3804	0.3804	0.0068 <sup>††</sup>	0.0122 <sup>†</sup>
<i>Bern_IPMI</i>	0.3804	0.4238	N/A	0.0122	0.0522	0.2661	0.0015 <sup>††</sup>	0.0015 <sup>††</sup>
<i>TU/e IMAG/e</i>	0.2036	0.0269	0.0122	N/A	0.6772	0.0342 <sup>†</sup>	0.1294	0.0342 <sup>†</sup>
<i>UPF_simbiosys</i>	0.4697	0.3804	0.0522	0.6772	N/A	0.2661	0.1294	0.0923
<i>NeuroMTL</i>	0.8501	0.3804	0.2661	0.0342 <sup>†</sup>	0.2661	N/A	0.0015 <sup>††</sup>	0.0049 <sup>††</sup>
<i>UPC_DLMI</i>	0.0068 <sup>††</sup>	0.0068 <sup>†</sup>	0.0015 <sup>††</sup>	0.1294	0.1294	0.0015 <sup>††</sup>	N/A	0.8501
<i>LRDE</i>	0.0640	0.0122 <sup>†</sup>	0.0015 <sup>††</sup>	0.0342 <sup>†</sup>	0.0923	0.0049 <sup>††</sup>	0.8501	N/A

Based above evaluation in terms of whole brain, small ROIs, and gyral curves, we can see none of method has achieved a significantly better performance over any other method. Especially, from the error maps in Fig. 12, these methods consistently have a poor performance along the cortical

regions. Therefore, these are still many spaces for improvement.

First, all methods directly apply well-established models (e.g., U-Nets) on the challenge, without considering any *prior* knowledge of infant brain images, e.g., cortical thickness is within a certain range. Especially, due to low contrast between

WM and GM in the 6-month infant brain images, WM voxels may be under/over segmented. Given a voxel with a resolution of  $1 \times 1 \times 1 \text{ mm}^3$ , although one voxel error will have a negligible impact on DICE or HD95, it will result in  $\pm 1 \text{ mm}$  estimation error of cortical thickness. Fig. 15 shows a segmentation result on a testing subject obtained by *MSL\_SKKU* [20]. Without anatomical guidance, there are many missing gyri in the reconstructed inner surface by *MSL\_SKKU* [20]. Consequently, the estimated cortical thickness is abnormally thicker. It is worth noting that this type of error should be paid more attention, especially for possible biomarker identification, since it might be difficult to accurately characterize brain developmental attributes, such as cortical thickness, gyrification, and convexity. For example, the cortical thickness for the zoomed regions (last column of Fig. 15) is abnormally larger than the ground truth.

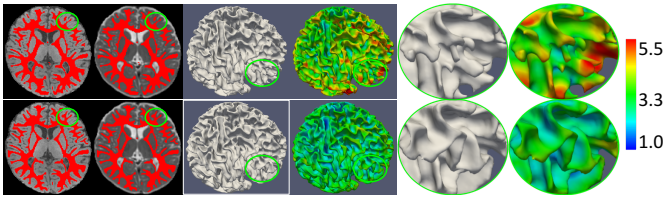


Fig. 15. Comparison with *MSL\_SKKU* [13] in 2017 MICCAI Grand Segmentation Challenge (iSeg-2017). The results by *MSL\_SKKU* [13] and manual segmentation are shown in the 1st and 2nd row, respectively. From left to right: segmentation overlaid on T1w and T2w images, inner surface, cortical thickness, zoomed views of inner surface and cortical thickness (in mm).

Second, all methods ignore a fact between CSF and GM is much higher than that between GM and WM. Therefore, it might be reasonable to identify CSF firstly from infant brain images to reconstruct the outer cortical surface and use it as a guidance to estimate the inner cortical surface since cortical thickness is within a certain range. Preliminary work on 6-month infant subjects with risk of autism demonstrates its effectiveness [53].

Third, among 13 testing subjects, we find all methods consistently performed badly on the 20-th testing subject. One representative slice is shown in Fig. 16. It can be seen that the image was with severe motion artifacts. By contrast the other subjects have small or no motion effect. Another possible reason could be from the different scan pose for this subject. Therefore, the model that is robust to the motion/scan pose is highly desired, since these artifacts are inevitable during the scanning. Possible solution is to augment the training subjects by rotation, flipping, as well as enforcing motion artifacts.



Fig. 16. The 20-th testing subject with motions and different scan pose.

Fourth, Table. V lists key highlights, time cost for training and testing, limitations, etc., for top-8 methods. All the methods typically randomly or gridly selected samples (2D/3D patches) from the training images, without realizing the importance of sample selection. For example, in conventional machine learning algorithms, ad-boosting is an effective strategy to learning features from these error-prone regions to improve the performance [54]. Similarly, the average error map shown in Fig. 12 could potentially provide guidance for samples selection. For example, by selecting more training samples from these error-prone regions, the performance of these deep learning algorithms could be further improved. In addition, from Table. V, we can see the patch sizes of deep learning models are varying from  $24 \times 24 \times 24$  to  $80 \times 80 \times 80$ , which could be further optimized.

There are also limitations for iSeg-2017. For example, the numbers of training subjects and testing subjects are limited. Another limitation is low image resolution, especially for T2w images with  $1.25 \times 1.25 \times 1.95 \text{ mm}^3$ . While currently, T1w and T2w images are usually acquired with  $1.0 \times 1.0 \times 1.0 \text{ mm}^3$ , or even  $0.8 \times 0.8 \times 0.8 \text{ mm}^3$  in BCP imaging protocol [1]. These limitations will be alleviated by including more subjects from BCP in the 2019 iSeg Segmentation Grand Challenge (<https://iseg2019.web.unc.edu>).

## VI. CONCLUSION

In this paper, we have reviewed and summarized 21 automatic segmentation methods participated in iSeg-2017. Especially, we have elaborated the details of the top 8 methods: including the pipeline, implementation, and source code. We further pointed out limitations and possible future directions. The iSeg-2017 website is still open and we hope our manual labels in iSeg-2017, this review article and source codes could greatly advance methodological development in the community.

## REFERENCES

- [1] B. R. Howell, M. A. Styner, W. Gao, P. T. Yap, L. Wang, K. Baluyot, *et al.*, "The UNC/UMN Baby Connectome Project (BCP): An overview of the study design and protocol development," *Neuroimage*, Mar 22 2018.
- [2] T. Paus, D. L. Collins, A. C. Evans, G. Leonard, B. Pike, and A. Zijdenbos, "Maturation of white matter in the human brain: a review of magnetic resonance studies," *Brain Res Bull*, vol. 54, pp. 255-66, Feb 2001.
- [3] G. Li, L. Wang, P.-T. Yap, F. Wang, Z. Wu, Y. Meng, *et al.*, "Computational neuroanatomy of baby brains: A review," *Neuroimage*, 2018/03/21/ 2018.

- [4] I. Isgum, M. J. N. L. Benders, B. Avants, M. J. Cardoso, S. J. Counsell, E. F. Gomez, *et al.*, "Evaluation of automatic neonatal brain segmentation algorithms: The NeoBrainS12 challenge," *Medical Image Analysis*, vol. 20, pp. 135-151, Feb 2015.
- [5] A. M. Mendrik, K. L. Vincken, H. J. Kuijff, M. Breeuwer, W. H. Bouvy, J. de Bresser, *et al.*, "MRBrainS Challenge: Online Evaluation Framework for Brain Image Segmentation in 3T MRI Scans," *Computational Intelligence and Neuroscience*, 2015.
- [6] S. Winzeck, A. Hakim, R. McKinley, J. A. A. D. S. R. Pinto, V. Alves, C. Silva, *et al.*, "ISLES 2016 and 2017-Benchmarking Ischemic Stroke Lesion Outcome Prediction Based on Multispectral MRI," *Frontiers in Neurology*, vol. 9, Sep 13 2018.



- [7] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, *et al.*, "The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS)," *IEEE Transactions on Medical Imaging*, vol. 34, pp. 1993-2024, Oct 2015.
- [8] W. L. Zhang, R. J. Li, H. T. Deng, L. Wang, W. L. Lin, S. W. Ji, *et al.*, "Deep convolutional neural networks for multi-modality isointense infant brain image segmentation," *Neuroimage*, vol. 108, pp. 214-224, Mar 2015.
- [9] L. Wang, F. Shi, Y. Gao, G. Li, J. H. Gilmore, W. Lin, *et al.*, "Integration of sparse multi-modality representation and anatomical constraint for isointense infant brain MR image segmentation," *Neuroimage*, vol. 89, pp. 152-64, Apr 1 2014.
- [10] D. Nie, L. Wang, E. Adeli, C. Lao, W. Lin, and D. Shen, "3-D fully convolutional networks for multimodal isointense infant brain image segmentation," *IEEE Transactions on Cybernetics*, 2018.
- [11] F. Shi, L. Wang, Y. K. Dai, J. H. Gilmore, W. L. Lin, and D. G. Shen, "LABEL: Pediatric brain extraction using learning-based meta-algorithm," *Neuroimage*, vol. 62, pp. 1975-1986, Sep 2012.
- [12] J. G. Sled, A. P. Zijdenbos, and A. C. Evans, "A nonparametric method for automatic correction of intensity nonuniformity in MRI data," *IEEE Transactions on Medical Imaging*, vol. 17, pp. 87-97, Feb 1998.
- [13] J. G. Chi, E. C. Dooling, and F. H. Gilles, "Gyral development of the human brain," *Ann Neurol*, vol. 1, pp. 86-93, Jan 1977.
- [14] E. Armstrong, A. Schleicher, H. Omran, M. Curtis, and K. Zilles, "The ontogeny of human gyrification," *Cereb Cortex*, vol. 5, pp. 56-63, Jan-Feb 1995.
- [15] J. Hill, D. Dierker, J. Neil, T. Inder, A. Knutsen, J. Harwell, *et al.*, "A surface-based analysis of hemispheric asymmetries and folding of cerebral cortex in term-born human infants," *J Neurosci*, vol. 30, pp. 2268-76, Feb 10 2010.
- [16] J. Dubois, M. Benders, A. Cachia, F. Lazeyras, R. Ha-Vinh Leuchter, S. V. Sizonenko, *et al.*, "Mapping the early cortical folding process in the preterm newborn brain," *Cereb Cortex*, vol. 18, pp. 1444-54, Jun 2008.
- [17] S. Abe, K. Takagi, T. Yamamoto, Y. Okuhata, and T. Kato, "Assessment of cortical gyrus and sulcus formation using MR images in normal fetuses," *Prenatal Diagnosis*, vol. 23, pp. 225-231, Mar 2003.
- [18] C. Lebel, L. Walker, A. Leemans, L. Phillips, and C. Beaulieu, "Microstructural maturation of the human brain from childhood to adulthood," *Neuroimage*, vol. 40, pp. 1044-1055, Apr 15 2008.
- [19] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3D U-Net: learning dense volumetric segmentation from sparse annotation," presented at the MICCAI, 2016.
- [20] T. D. Bui, J. Shin, and T. Moon, "3D Densely Convolutional Networks for Volumetric Segmentation," *arXiv:1709.03199*, 2017.
- [21] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely Connected Convolutional Networks," presented at the CVPR, 2017.
- [22] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [23] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," presented at the Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, 2011.
- [24] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, pp. 1929-1958, 2014.
- [25] J. Long, E. Shelhamer, and T. Darrell, "Fully Convolutional Networks for Semantic Segmentation," presented at the CVPR, 2015.
- [26] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [27] K. M. He, X. Y. Zhang, S. Q. Ren, and J. Sun, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification," presented at the ICCV, 2015.
- [28] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, *et al.*, "Caffe: Convolutional architecture for fast feature embedding," presented at the Proceedings of the 22nd ACM international conference on Multimedia, 2014.
- [29] D. C. Dolz, J. Wang, L. Yuan, J. Shen, D. Ayed, "Deep CNN ensembles and suggestive annotations for infant brain MRI segmentation," *arXiv preprint arXiv:1712.05319*, 2017.
- [30] J. Dolz, C. Desrosiers, and I. Ben Ayed, "3D fully convolutional networks for subcortical segmentation in MRI: A large-scale study," *Neuroimage*, vol. 170, pp. 456-470, Apr 15 2018.
- [31] G. Zeng and G. Zheng, "Multi-stream 3D FCN with multi-scale deep supervision for multi-modality isointense infant brain MR image segmentation," presented at the ISBI, 2018.
- [32] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," presented at the MICCAI, 2015.
- [33] R. Kimmel, N. Kiryati, and A. M. Bruckstein, "Sub-pixel distance maps and weighted distance transforms," *Journal of Mathematical Imaging and Vision*, vol. 6, pp. 223-233, Jun 1996.
- [34] P. Moeskops and J. P. Pluim, "Isointense infant brain MRI segmentation with a dilated convolutional neural network," *arXiv preprint arXiv:1708.02757*, 2017.
- [35] P. Moeskops, M. A. Viergever, A. M. Mendrik, L. S. de Vries, M. J. Benders, and I. Išgum, "Automatic segmentation of MR brain images with a convolutional neural network," *IEEE Transactions on Medical Imaging*, vol. 35, pp. 1252-1261, 2016.
- [36] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [37] J. M. Wolterink, T. Leiner, M. A. Viergever, and I. Išgum, "Dilated convolutional neural networks for cardiovascular MR segmentation in congenital heart disease," in *Reconstruction, Segmentation, and Analysis of Medical Images*, ed: Springer, 2016, pp. 95-102.
- [38] G. Sanroma, O. M. Benkarim, G. Piella, and M. A. G. Ballester, "Building an Ensemble of Complementary Segmentation Methods by Exploiting Probabilistic Estimates," *Machine Learning in Medical Imaging, Mlmi 2016*, vol. 10019, pp. 27-35, 2016.
- [39] H. Z. Wang, J. W. Suh, S. R. Das, J. B. Pluta, C. Craige, and P. A. Yushkevich, "Multi-Atlas Segmentation with Joint Label Fusion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, pp. 611-623, Mar 2013.
- [40] Y. F. Hao, T. Y. Wang, X. Q. Zhang, Y. Y. Duan, C. S. Yu, T. Z. Jiang, *et al.*, "Local Label Learning (LLL) for Subcortical Structure Segmentation: Application to Hippocampus Segmentation," *Human Brain Mapping*, vol. 35, pp. 2674-2697, Jun 2014.
- [41] P. Moeskops, M. J. N. L. Benders, S. Chita, K. J. Kersbergen, F. Groenendaal, L. S. de Vries, *et al.*, "Automatic segmentation of MR brain images of preterm infants using supervised classification," *Neuroimage*, vol. 118, pp. 628-641, Sep 2015.
- [42] B. B. Avants, N. J. Tustison, G. Song, P. A. Cook, A. Klein, and J. C. Gee, "A reproducible evaluation of ANTs similarity metric performance in brain image registration," *Neuroimage*, vol. 54, pp. 2033-2044, Feb 1 2011.
- [43] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [44] V. Fonov, A. Doyle, A. C. Evans, and D. L. Collins, "NeuroMTL iSEG challenge methods," *bioRxiv*, 2018.
- [45] H. C. Hazlett, H. B. Gu, B. C. Munsell, S. H. Kim, M. Styner, J. J. Wolff, *et al.*, "Early brain development in infants at high risk for autism spectrum disorder," *Nature*, vol. 542, pp. 348-351, Feb 16 2017.
- [46] V. Fonov, A. C. Evans, K. Botteron, C. R. Almli, R. C. McKinstry, D. L. Collins, *et al.*, "Unbiased average age-appropriate atlases for pediatric studies," *Neuroimage*, vol. 54, pp. 313-327, Jan 1 2011.
- [47] N. N. Fausto Milletari, Seyed-Ahmad Ahmadi, "V-Net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation," *arXiv:1606.04797*, 2016.
- [48] Y. C. Xu, T. Geraud, and I. Bloch, "From Neonatal to Adult Brain MR Image Segmentation in a Few Seconds Using 3D-Like Fully Convolutional Network and Transfer Learning," presented at the ICIP, 2017.
- [49] I. S. Gousias, D. Rueckert, R. A. Heckemann, L. E. Dyet, J. P. Boardman, A. D. Edwards, *et al.*, "Automatic segmentation of brain MRIs of 2-year-olds into 83 regions of interest," *Neuroimage*, vol. 40, pp. 672-684, Apr 1 2008.
- [50] B. Fischl, "FreeSurfer," *Neuroimage*, vol. 62, pp. 774-81, Aug 15 2012.
- [51] B. B. Avants, C. L. Epstein, M. Grossman, and J. C. Gee, "Symmetric diffeomorphic image registration with cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain," *Med Image Anal*, vol. 12, pp. 26-41, Feb 2008.
- [52] Y. Zhang, F. Shi, G. Wu, L. Wang, P.-T. Yap, and D. Shen, "Consistent spatial-temporal longitudinal atlas construction for developing infant brains," *IEEE transactions on medical imaging*, vol. 35, pp. 2568-2577, 2016.
- [53] L. Wang, G. Li, F. Shi, X. Cao, C. Lian, D. Nie, *et al.*, "Volume-based analysis of 6-month-old infant brain MRI for autism biomarker identification and early diagnosis," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2018, pp. 411-419.
- [54] S. Shalev-Shwartz, "Selfieboost: A boosting algorithm for deep learning," *arXiv preprint arXiv:1411.3436*, 2014.

Table IV. AVERAGE DICE SCORES FOR 80 ROIS BY ALL THE COMPETING METHODS. L: LEFT; R: RIGHT.

ROI\Method	MSL_SKKU	LIVIA	Bern_IPMI	TU/e IMAG/e	UPF_simbiosys	NeuroMTL	UPC_DLMI	LRDE	Average
L hippocampus	0.8932	0.8851	0.8944	0.8606	0.8838	0.8774	0.8676	0.8374	0.8749
R hippocampus	0.8988	0.9029	0.9063	0.8726	0.9009	0.8997	0.9027	0.8752	0.8949
L amygdala	0.9209	0.8937	0.9304	0.9036	0.8892	0.8679	0.8982	0.8714	0.8969
R amygdala	0.9126	0.8878	0.9064	0.8911	0.8946	0.9005	0.8978	0.8818	0.8966
L anterior temporal lobe, medial part	0.9062	0.8981	0.8902	0.8932	0.8739	0.8293	0.8871	0.8637	0.8802
R anterior temporal lobe, medial part	0.8993	0.8902	0.882	0.8801	0.8765	0.8105	0.8795	0.8517	0.8712
L anterior temporal lobe, lateral part	0.9202	0.9211	0.9036	0.9072	0.8916	0.8551	0.8925	0.8827	0.8968
R anterior temporal lobe, lateral part	0.9159	0.9207	0.8976	0.9012	0.8960	0.8431	0.8912	0.8820	0.8935
L parahippocampal and ambient gyri	0.8838	0.8696	0.8755	0.8634	0.8672	0.8583	0.8427	0.8281	0.8611
R parahippocampal and ambient gyri	0.8870	0.8767	0.8816	0.8596	0.8736	0.8617	0.8415	0.8330	0.8643
L superior temporal gyrus, posterior part	0.9253	0.9177	0.9196	0.9060	0.9008	0.9085	0.9073	0.8833	0.9086
R superior temporal gyrus, posterior part	0.9217	0.9161	0.9188	0.9036	0.9008	0.9128	0.9045	0.8835	0.9077
L middle and inferior temporal gyrus	0.9100	0.8988	0.8874	0.8899	0.8775	0.8629	0.8777	0.8610	0.8832
R middle and inferior temporal gyrus	0.9015	0.9001	0.8817	0.8904	0.8732	0.8653	0.8756	0.8554	0.8804
L fusiform gyrus	0.9034	0.8969	0.8693	0.8859	0.8766	0.8028	0.8721	0.8539	0.8701
R fusiform gyrus	0.8865	0.8796	0.8589	0.8710	0.8677	0.7885	0.8571	0.8333	0.8553
L insula	0.8954	0.8880	0.8891	0.8825	0.8823	0.8913	0.8776	0.8740	0.8850
R insula	0.9040	0.8891	0.8983	0.8857	0.8890	0.8943	0.8932	0.8755	0.8911
L lateral remainder of occipital lobe	0.9000	0.9040	0.8908	0.8883	0.8837	0.8915	0.8979	0.8617	0.8897
R lateral remainder of occipital lobe	0.9019	0.8986	0.8883	0.8917	0.8828	0.8829	0.8915	0.8601	0.8872
L cingulate gyrus, anterior part	0.9255	0.9188	0.9148	0.9000	0.9068	0.9135	0.9040	0.8925	0.9095
R cingulate gyrus, anterior part	0.9254	0.9203	0.9193	0.9083	0.9070	0.9121	0.9156	0.8911	0.9124
L cingulate gyrus, posterior part	0.9129	0.9141	0.9094	0.8912	0.9036	0.9094	0.8976	0.8819	0.9025
R cingulate gyrus, posterior part	0.9162	0.9172	0.9002	0.8972	0.9013	0.8974	0.9028	0.8757	0.9010
L middle frontal gyrus	0.9349	0.9308	0.9235	0.9273	0.9136	0.9208	0.9228	0.9077	0.9227
R middle frontal gyrus	0.9317	0.9290	0.9193	0.9222	0.9116	0.9158	0.9173	0.9030	0.9187
L Posterior temporal lobe	0.9123	0.9103	0.8945	0.9002	0.8939	0.8998	0.8985	0.8767	0.8983
R posterior temporal lobe	0.9176	0.9138	0.9016	0.9013	0.8974	0.8916	0.9015	0.8797	0.9006
L inferiolateral remainder of parietal lobe	0.9258	0.9269	0.9111	0.9181	0.8970	0.9151	0.9164	0.8942	0.9131
R inferiolateral remainder of parietal lobe	0.9276	0.9219	0.9124	0.9131	0.8943	0.9086	0.9103	0.8902	0.9098
L caudate nucleus	0.9282	0.9310	0.9360	0.8972	0.9374	0.9066	0.9219	0.9094	0.9210
R caudate nucleus	0.9261	0.9284	0.9311	0.8739	0.9335	0.9096	0.9276	0.9119	0.9178
L nucleus accumbens	0.8469	0.8404	0.8682	0.8491	0.9228	0.7665	0.8830	0.8551	0.8540
R nucleus accumbens	0.8648	0.8949	0.8588	0.8549	0.9409	0.7330	0.9058	0.9080	0.8701
L putamen	0.9350	0.9410	0.9298	0.9299	0.9501	0.9024	0.9480	0.9170	0.9317
R putamen	0.9330	0.9499	0.9420	0.9520	0.9568	0.9138	0.9453	0.9382	0.9414
L thalamus	0.9249	0.9303	0.9418	0.9081	0.9438	0.9345	0.9337	0.9236	0.9301
R thalamus	0.9241	0.9270	0.9457	0.9203	0.9417	0.9400	0.9352	0.9294	0.9329
L pallidum	0.7410	0.7529	0.7577	0.7343	0.7756	0.7429	0.7631	0.7193	0.7484
R pallidum	0.7354	0.7679	0.7658	0.7671	0.7440	0.7657	0.7432	0.7639	0.7566
Corpus callosum	0.9469	0.9385	0.9412	0.9378	0.9364	0.9377	0.9394	0.9231	0.9376
L lateral ventricle (excluding temporal horn)	0.9427	0.9409	0.9437	0.9212	0.9263	0.9319	0.9286	0.9204	0.9320
R lateral ventricle (excluding temporal horn)	0.9216	0.9243	0.9234	0.9004	0.8987	0.9099	0.8925	0.8912	0.9078

L lateral ventricle, temporal horn	0.8077	0.8275	0.8014	0.7836	0.7949	0.7800	0.7703	0.7552	0.7901
R lateral ventricle, temporal horn	0.8155	0.8576	0.8373	0.8257	0.8197	0.7756	0.8005	0.7724	0.8130
Third ventricle	0.9587	0.9612	0.9629	0.9449	0.9555	0.9545	0.9546	0.9391	0.9539
L precentral gyrus	0.9342	0.9273	0.9316	0.9186	0.9164	0.9234	0.9251	0.9022	0.9224
R precentral gyrus	0.9331	0.9259	0.929	0.9166	0.9154	0.9165	0.9204	0.9013	0.9198
L straight gyrus	0.8534	0.8375	0.8309	0.8411	0.8252	0.8398	0.8389	0.7961	0.8329
R straight gyrus	0.8672	0.8703	0.8552	0.8715	0.8308	0.8417	0.8594	0.8131	0.8512
L anterior orbital gyrus	0.9023	0.8898	0.8765	0.8923	0.8768	0.8733	0.8762	0.8519	0.8799
R anterior orbital gyrus	0.8975	0.8867	0.8746	0.8864	0.8778	0.8720	0.8738	0.8485	0.8772
L inferior frontal gyrus	0.9251	0.919	0.9173	0.9108	0.8985	0.9109	0.9059	0.8895	0.9096
R inferior frontal gyrus	0.9251	0.9185	0.9167	0.9063	0.9001	0.9030	0.9048	0.8907	0.9082
L superior frontal gyrus	0.9281	0.9247	0.9193	0.9157	0.9118	0.9171	0.913	0.8982	0.916
R superior frontal gyrus	0.9271	0.9231	0.9203	0.9153	0.9090	0.9155	0.9143	0.8979	0.9153
L postcentral gyrus	0.9252	0.9175	0.9191	0.9037	0.9045	0.9116	0.9081	0.8885	0.9098
R postcentral gyrus	0.9255	0.9176	0.9214	0.9068	0.9040	0.9096	0.9084	0.8852	0.9098
L superior parietal gyrus	0.9213	0.9155	0.9128	0.9093	0.8968	0.9078	0.9107	0.8808	0.9069
R superior parietal gyrus	0.9173	0.9145	0.9100	0.9049	0.8946	0.9063	0.9054	0.8794	0.9041
L lingual gyrus	0.8864	0.8800	0.8650	0.8606	0.8664	0.8598	0.8595	0.8260	0.8630
R lingual gyrus	0.8785	0.8717	0.8553	0.8533	0.8563	0.8513	0.8515	0.8104	0.8535
L cuneus	0.8922	0.8839	0.8732	0.8685	0.8708	0.8699	0.8668	0.8306	0.8695
R cuneus	0.8848	0.8783	0.8770	0.8614	0.8654	0.8659	0.8641	0.8258	0.8653
L medial orbital gyrus	0.8959	0.8809	0.8685	0.8771	0.8828	0.8593	0.8645	0.8466	0.8720
R medial orbital gyrus	0.8863	0.8809	0.866	0.8768	0.8798	0.8632	0.8668	0.8376	0.8697
L lateral orbital gyrus	0.8992	0.8883	0.8863	0.8924	0.8744	0.8766	0.8649	0.8596	0.8802
R lateral orbital gyrus	0.8963	0.8918	0.8807	0.9004	0.8791	0.8777	0.8645	0.8527	0.8804
L posterior orbital gyrus	0.8925	0.8804	0.8789	0.8835	0.8717	0.8661	0.8583	0.8490	0.8726
R posterior orbital gyrus	0.9035	0.8970	0.8935	0.8906	0.8830	0.8721	0.8684	0.8569	0.8831
L substantia nigra	0.9295	0.9213	0.9232	0.9193	0.9157	0.9292	0.8632	0.8809	0.9103
R substantia nigra	0.9225	0.9150	0.9208	0.9121	0.9178	0.9076	0.8735	0.8860	0.9069
L subgenual frontal cortex	0.8700	0.8452	0.8531	0.8482	0.8575	0.8648	0.8550	0.8454	0.8549
R subgenual frontal cortex	0.8739	0.8667	0.8476	0.8680	0.8486	0.8694	0.8617	0.8386	0.8593
L subcallosal area	0.9012	0.8760	0.8731	0.8291	0.8758	0.8669	0.8845	0.8738	0.8726
R subcallosal area	0.8376	0.8551	0.8626	0.8200	0.8437	0.8576	0.8563	0.8465	0.8474
L pre-subgenual frontal cortex	0.8994	0.8974	0.8967	0.8885	0.8850	0.9062	0.8777	0.8870	0.8922
R pre-subgenual frontal cortex	0.8801	0.8824	0.8520	0.8772	0.8544	0.8757	0.8815	0.8683	0.8715
L superior temporal gyrus, anterior part	0.9127	0.9024	0.9106	0.8908	0.8969	0.8991	0.8878	0.8787	0.8974
R superior temporal gyrus, anterior part	0.9126	0.9018	0.9028	0.8903	0.8807	0.8874	0.8881	0.8697	0.8917



TABLE V. PARAMETERS OF TOP-8 TEAMS IN TERMS OF ARCHITERTURE, TOOL, HIGHLIGHT, ETC.

Team	Architecture	Tool	Key highlights	Augmentation	Pretrained?	Training Loss	Memory (training/testing)	Time (training/testing)	2D/3D Patch size (training/testing)	Patch selection	Limitation
<i>MSL_SKKU</i>	DenseNet	Caffe	Skip-connections	No	No	Cross-entropy	12 G/2G	2 days/300 seconds	3D (64×64×64/64×64×64)	Random	No prior/No augmentation/No ad-boosting
<i>LIVIA</i>	Semi-dense or quasi-dense FCN	Theano	Quasi-dense architecture; Ensemble	No	No	Cross-entropy	6 G/2G	2 days/~10 seconds	3D (27×27×27/35×35×35)	Random	No prior/No augmentation/No ad-boosting
<i>Bern_IPMI</i>	Stacked U-Nets	Tensorflow	Multi-Scale, Two stages	Rotation and flip, distance maps	Yes	Cross-entropy	11 G/11G	8 hours/8 seconds	3D (64×64×64/64×64×64)	Random	No prior/No ad-boosting
<i>TU/e IMAG/e</i>	Dilated convolutional neural network	Lasagne + Theano	2.5D dilated CNN combined with 3D CNN	No	No	Cross-entropy	12G/2G	1 day/60 seconds	2.5D+3D (67×67+25×25×25/67×67+25×25×25)	Random	No prior/No augmentation/No ad-boosting
<i>UPF_simbiosys</i>	Cascading (multi-atlas label fusion + SVM)	SVM+ ANTs	Spatial priors + multi-scale+ SVM	Registration priors from multi-atlas label fusion	N/A	N/A	4G/4G	1 hour/30 mins	3D (5×5×5/5×5×5)	Random	High computational time required at testing
<i>NeuroMTL</i>	U-Net	Torch	Pre-training from another age-matched dataset	No	Yes	Cross-entropy	12G/2.2G	11 hours/8 seconds	3D (80×80×80/80×80×80)	Grid	No prior/No ad-boosting
<i>UPC_DLMI</i>	U-Net	Tensorflow + Keras	Augmented path for learning high resolution features	Flipping	No	Weighted cross-entropy	12G/8G	1.5 days/7 seconds	3D (64*64*64/64*64*64)	Random	No prior/No ad-boosting
<i>LRDE</i>	VGG-like+FCN	Caffe	VGG for brain segmentation; Building 2D color image from 3D MRI Volume	Rotations, translation, scaling	Yes	Cross-entropy	8G/1G	4 hours/1.8 seconds	3D (x*y*3/x*y*3, 3 channels)	Random	Discontinuities between slices due to only 2D patch involved/No prior/No ad-boosting